

# Détection de fraudes et loi de probabilité de Newcomb-Benford

Projet Master 1

FERNANDEZ Christelle      PONCHEELE Clément      EL KAÏM Laura  
Encadré par M.DUCHARME

24 *mai* 2021

# Remerciements

Nous souhaitons remercier la faculté des Sciences de Montpellier pour les Master MIND (Mathématiques de l'Information et de la Décision) et Biostatistiques et plus particulièrement Monsieur **Ducharme** pour nous avoir permis de réaliser ce sujet.

Lors de ce projet, nous avons pu affiner notre travail en équipe et notre autonomie, consolider nos acquis et ce rapport signe l'aboutissement de notre première année de Master.

Nous remercions également nos proches qui nous ont soutenu dans l'élaboration de notre projet et remercions notamment les participants à notre expérience.

Remerciements spéciaux à nos relecteurs et correcteurs qui ont contribué au bon déroulement du rapport.

# Résumé

Dans différents cadres, la fraude est existante. Une façon répandue pour commettre une fraude est de modifier des chiffres de données de la manière dont le désire l'escroc et notamment en modifiant le premier chiffre significatif. Ce que nous étudierons plus particulièrement ici.

Pour détecter ces fraudes, nous pouvons utiliser la loi de Newcomb-Benford sur des échantillons de données.

L'objectif étant d'exposer l'émergence de la loi de Newcomb-Benford dans des données réelles, d'étudier différents jeux de données comme ici la fiscalité italienne et à l'aide de divers tests statistiques suspecter ou non une fraude, par le biais de la modification du premier chiffre significatif des nombres. Nous pourrions utiliser différents tests et voir si ceux-ci vont plutôt dans le même sens ou si certains se contredisent. Nous aborderons les tests dits classiques d'adéquation à la loi de Newcomb-Benford, ainsi que les tests lisses mis en place par Monsieur Ducharme et ses collaborateurs.

Pour répondre à nos différentes problématiques, nous allons effectuer des expériences d'abord visuelles, en extrayant des données d'un journal, d'un magazine ainsi que d'autres données réelles puis nous appliquerons différents tests sur ces jeux de données à l'aide du logiciel R et des packages **BenfordTests** et **BENFORDSMOOTHTEST**. Nous utiliserons également ces packages sur un jeu de données fiscales italiennes.

Les réponses récoltées dans le premier temps, nous montrerons que la loi de Newcomb-Benford n'apparaît pas partout notamment sur les données influencées par la pensée de l'homme, les données dites "non-naturelles". Nous observons également qu'une inspection visuelle n'est pas suffisante pour suspecter ou non une fraude. Puis dans un second temps, les tests réalisés nous permettrons de voir qu'il est parfois difficile de suspecter une fraude, au risque de se tromper.

**Mots-clés** : Loi de Newcomb-Benford, 1<sup>er</sup> chiffre significatif, Test d'hypothèses, Hypothèse nulle/alternative, Risque d'erreur, Test d'adéquation, Test du khi-deux, Test lisse, Test de Freedman, p-value.

# Table des matières

<b>Remerciements</b>	<b>i</b>
<b>Résumé</b>	<b>i</b>
<b>Introduction</b>	<b>1</b>
<b>Génèse de la loi de Newcomb-Benford</b>	<b>2</b>
<b>Expérimentation sur différents jeux de données</b>	<b>4</b>
La suite de Fibonacci . . . . .	4
Nombres extraits d'un magazine et d'un journal . . . . .	5
Population des villes de France . . . . .	6
Passage journalier de vélos dans l'allée Beracasa à Montpellier . . . . .	7
Nombres générés par les humains . . . . .	7
<b>Tests</b>	<b>9</b>
Généralités sur les tests . . . . .	9
Hypothèse nulle et hypothèse alternative . . . . .	10
Les risques d'erreurs . . . . .	10
Test du Khi-Deux . . . . .	12
Test de Freedman-Watson . . . . .	13
Tests lisses pour la Loi de Newcomb-Benford . . . . .	14
<b>Application des tests</b>	<b>16</b>
Application sur nos jeux de données . . . . .	16
Application a des données fiscales italiennes . . . . .	18
CAS COVID . . . . .	23
<b>Conclusion</b>	<b>24</b>
<b>Bibliographie</b>	<b>25</b>

## Table des figures

1	Figure 1 : Table logarithmique recueillie par M. Ducharme présentant des marques d'usure plus importantes sur les premières pages . . . . .	2
2	Figure 2 : Histogramme de la répartition du 1er chiffre significatif de la suite de Fibonacci en comparaison avec la loi de Benford . . . . .	4
3	Figure 3 : Histogramme de la répartition du 1er chiffre significatif des prix du magazine AMPM en comparaison avec la loi de Benford . . . . .	5
4	Figure 4 : Histogramme de la répartition du 1er chiffre significatif des nombres du journal LES ECHOS en comparaison avec la loi de Benford . . . . .	5
5	Figure 5 : Histogramme de la répartition du 1er chiffre significatif de la population française en comparaison avec la loi de Benford . . . . .	6
6	Figure 6 : Histogramme de la répartition du 1er chiffre significatif du passage journalier de vélos en comparaison avec la loi de Benford . . . . .	7
7	Figure 7 : Comparaison des histogrammes de la répartition du 1er chiffre significatif de l'expérience de Hill avec la loi de Benford et la loi uniforme puis de notre expérience avec ces deux mêmes lois . . . . .	8
8	Figure 8 : Densité de la loi du Khi-Deux en fonction du nombre de degrés de liberté . .	12

## Liste des tableaux

1	Tableau 1 : Répartition du premier chiffre significatif selon la loi de Newcomb-Benford .	3
2	Tableau 2 : Règle de décision et risques d'erreurs . . . . .	11
3	Tableau 3 : Tests sur la région d'Abruzzo . . . . .	18
4	Tableau 4 : Tests sur la région de Basilicata . . . . .	18
5	Tableau 5 : Tests sur la région de Calabria . . . . .	18
6	Tableau 6 : Tests sur la région de la Lazio . . . . .	19
7	Tableau 7 : Tests sur la région de Friuli-Venezia-Giulia . . . . .	19
8	Tableau 8 : Tests sur la région de Lombardia . . . . .	19
9	Tableau 9 : Tests sur la région de Piemonte . . . . .	19
10	Tableau 10 : Tests sur la région de Puglia . . . . .	20
11	Tableau 11 : Tests sur la région de Trentino-alto . . . . .	20
12	Tableau 12 : Tests sur la région de Valle D'aosta . . . . .	20
13	Tableau 13 : Tests sur la région de la Sicilia . . . . .	20
14	Tableau 15 : Tests sur la région de Veneto . . . . .	21
15	Tableau 15 : Tests sur la région de Campania . . . . .	21
16	Tableau 16 : Tests sur la région d'Emilia-romagna . . . . .	21
17	Tableau 17 : Tests sur la région de la Liguria . . . . .	21
18	Tableau 18 : Tests sur la région de Marche . . . . .	22
19	Tableau 19 : Tests sur la région de Molise . . . . .	22
20	Tableau 20 : Tests sur la région de la Sardegna . . . . .	22
21	Tableau 21 : Tests sur la région de Toscana . . . . .	22
22	Tableau 22 : Tests sur la région de Umbria . . . . .	23

# Introduction

La fraude est une pratique répandue dans de nombreux domaines comme la finance, le secteur social ou médical. Il peut être tentant pour un être humain ou une société de tricher si cela peut impliquer pour lui une position plus confortable dans la société, telle qu'une réduction de charges, ou même un avantage sur un de ses concurrents. Il semblerait donc logique que des personnes cherchent à déceler ces fraudes.

Les données transmises par un individu ou un organisme peuvent faire l'objet de modifications, c'est de ce type de fraudes auquel nous nous intéresserons ici, et plus particulièrement la modification du premier chiffre significatif (le premier chiffre d'un nombre qui n'est pas un zéro) de nombres pris dans un certain ensemble de données.

De telles modifications entraînent un changement de la répartition des chiffres présents naturellement<sup>1</sup>. Si nous connaissons la répartition des chiffres présentés dans un ensemble de données arbitraires, il est donc techniquement possible de savoir si un nombre a été modifié ou non.

Il nous vient donc les questions suivantes : *Qu'elle est cette répartition ? Est-il possible de la connaître et si oui, dans quels cas ?*

De manière intuitive, nous pourrions penser que les nombres sont répartis de manière uniforme. Qu'en est-il vraiment ?

La première partie de notre projet consistera à **répondre à ces questions**, nous nous appuierons sur les travaux de Simon Newcomb et Frank Benford, qui ont théorisé la **loi de Newcomb-Benford**, plus communément appelée loi de Benford. Cette loi nous dit que, dans une liste de données dites naturelles, la probabilité d'avoir le chiffre  $i$  comme premier chiffre significatif est de  $\log_{10}(1 + \frac{1}{i})$ .

Par exemple, le chiffre 1 en tant que premier chiffre significatif serait présent à hauteur de 30% alors que le 9 à seulement 4,6%.

Dans la suite **nous mettrons en œuvre une série d'expérimentations** pour constater ou non la véracité de cette loi, pour ce faire dans un premier temps, nous récolterons des nombres pris dans des milieux censés satisfaire la loi de Newcomb-Benford et observerons la répartition du premier chiffre significatif. Puis nous répliquerons une version simplifiée de l'expérience de Hill (1988), qui consiste à observer la répartition du premier chiffre significatif d'une liste de nombre donnée au hasard par des êtres humains, en l'occurrence ses élèves.

Cette expérience est à la base des méthodes de détection de fraudes par la loi de Newcomb-Benford. Si un fraudeur modifie un jeu de données, ce jeu est donc influencé par la pensée humaine, il ne suit donc plus la loi de Newcomb-Benford. Pour détecter la fraude, il suffit donc de comparer les premiers chiffres significatifs. Cependant, ces comparaisons doivent se faire de manière rigoureuse et scientifique. Pour cela il existe des tests statistiques, dont le plus connu, le test du  $\chi^2$ , ou bien celui de Ducharme et collab. (2020). # ou bien des tests lisses réalisés par Ducharme et collab ?

Il nous vient donc les questions suivantes : *Ces tests sont-ils fiables ? Existe-t-il un test significativement meilleur que les autres ? Vont-ils dans le même sens ? Et sinon que faire ?*

La réponse à ces questions constituera donc la deuxième partie de ce projet, pour ce faire nous mettrons en œuvre différents tests sur des jeux de données comme la fiscalité italienne.

---

1. Les données dites naturelles sont celles qui n'ont pas été influencé par la pensée de l'homme.

# Génèse de la loi de Newcomb-Benford

Il serait tentant de penser que les nombres sont répartis de manière uniforme, cela viendrait du biais d'équiprobabilité<sup>2</sup>. Ce dernier consiste à “penser qu'en l'absence d'information, tous les cas ont la même probabilité de se produire et que le hasard implique nécessairement l'uniformité”.

Néanmoins, cette hypothèse sera contredite une première fois par l'astronome, mathématicien, économiste et statisticien canadien Simon Newcomb. Ce dernier fournira en 1881 une première approche au principe statistique, qui se fera injustement appeler *Loi de Benford*. Celui-ci remarquera que les premières pages des tables logarithmiques sont plus utilisées que les pages suivantes (cf. Figure 1). Il publiera sa découverte dans un article de l’*“American Journal of Mathematics”*.

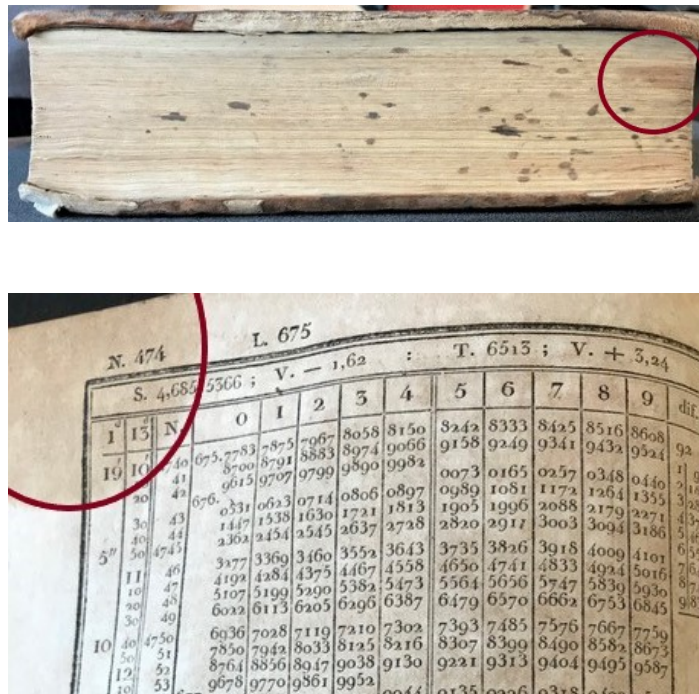


Figure 1 : Table logarithmique recueillie par M. Ducharme présentant des marques d'usure plus importantes sur les premières pages

Cette découverte mise de côté pendant plusieurs années, ce n'est qu'en 1938 que l'ingénieur et physicien américain Frank Benford arrivera au même résultat après avoir répertorié des dizaines de milliers de données. Celui-ci pensera être le premier à l'initiative de cette loi, et c'est pour cette raison que la *loi de Newcomb-Benford* se fera plus généralement appelée *loi de Benford*.

Cette loi nous dit que, dans une liste de données arbitraires, la probabilité d'avoir le chiffre  $i$  comme premier chiffre significatif est de  $\log_{10}(1 + \frac{1}{i})$ .

Nous retrouvons cette loi dans énormément de domaines comme les mathématiques, l'environnement, la finance, la physique, etc, plus précisément sur des données telles que la longueur des fleuves, la

2. Défini en 1985 par Marie-Paule Lecoutre (*Source*).

<b>PCS</b>	1	2	3	4	5	6	7	8	9
<b>Benford</b>	0.301	0.176	0.125	0.097	0.079	0.067	0.058	0.051	0.046

Tableau 1 : Répartition du premier chiffre significatif selon la loi de Newcomb-Benford

population des villes dans un pays, des déclarations de revenus, etc.

Notons cependant qu'il existe des cas où les données ne suivent pas cette loi, notamment des données dites non naturelles qui seraient influencé par la pensée humaine (nombres premiers, nombres générés par des humains, etc).



# Expérimentation sur différents jeux de données

Après avoir pris connaissance de la **loi de Newcomb-Benford**, il serait intéressant de la mettre en pratique sur différents jeux de données.

## La suite de Fibonacci

Intéressons-nous dans un premier temps à la suite de Fibonacci.

Cette suite est une suite d'entiers dans laquelle chaque terme est la somme des deux termes qui le précèdent. Sa formulation est la suivante :

$$F_0 = 0, F_1 = 1, \text{ et } \forall n \geq 2, F_n = F_{n-1} + F_{n-2}.$$

Nous commençons par recueillir les 1000 premiers termes de la suite de Fibonacci, pour extraire le premier chiffre significatif de chacun de ces nombres.

Par la suite nous calculons la répartition de chaque chiffre significatif et obtenons l'histogramme suivant, où  $n$  ainsi que pour l'ensemble des histogrammes suivants, représente le nombre de données de l'échantillon :

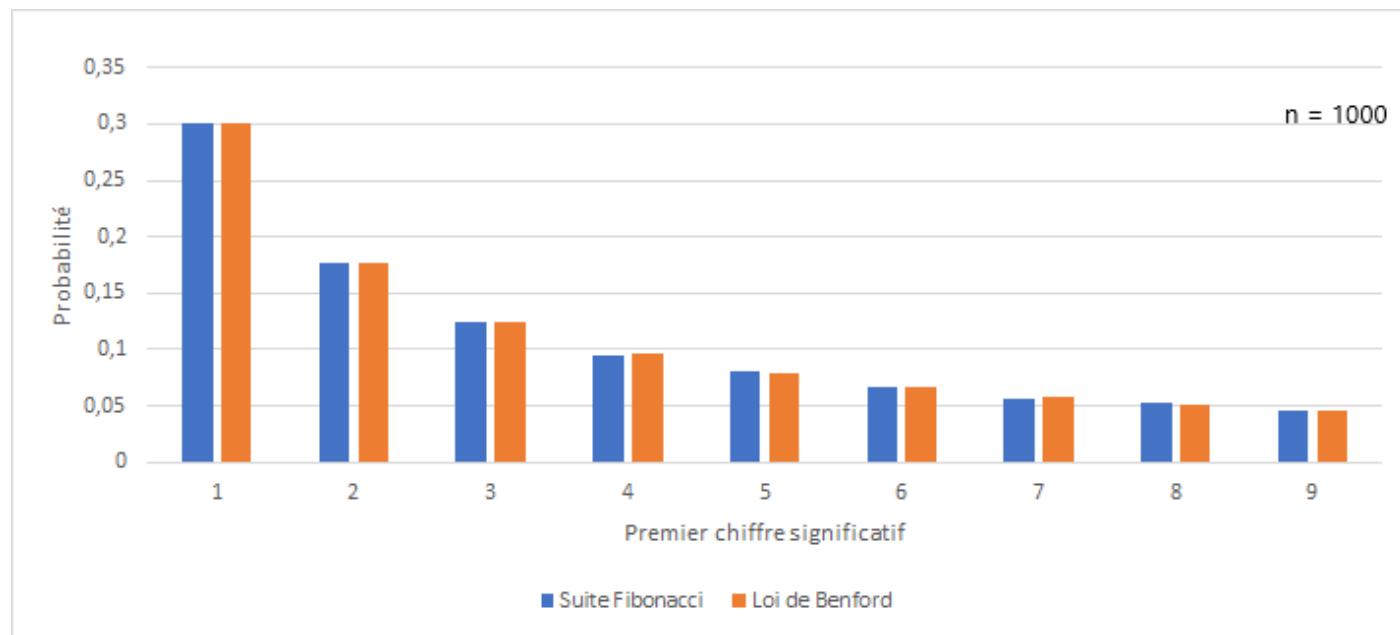


Figure 2 : Histogramme de la répartition du 1er chiffre significatif de la suite de Fibonacci en comparaison avec la loi de Benford

Visuellement, il semblerait que la répartition des chiffres significatifs des 1000 premiers nombres de la suite de Fibonacci suive la **loi de Newcomb-Benford**.

## Nombres extraits d'un magazine et d'un journal

Dans un second temps, nous relevons les prix présents dans un magazine de mobilier de la marque *AMPM*, ainsi que tous les nombres répertoriés dans un journal *Les ECHOS*. Nous récoltons environ 300 nombres par magazine et, de la même façon qu'énoncé précédemment, calculons la répartition des chiffres significatifs de ces nombres.

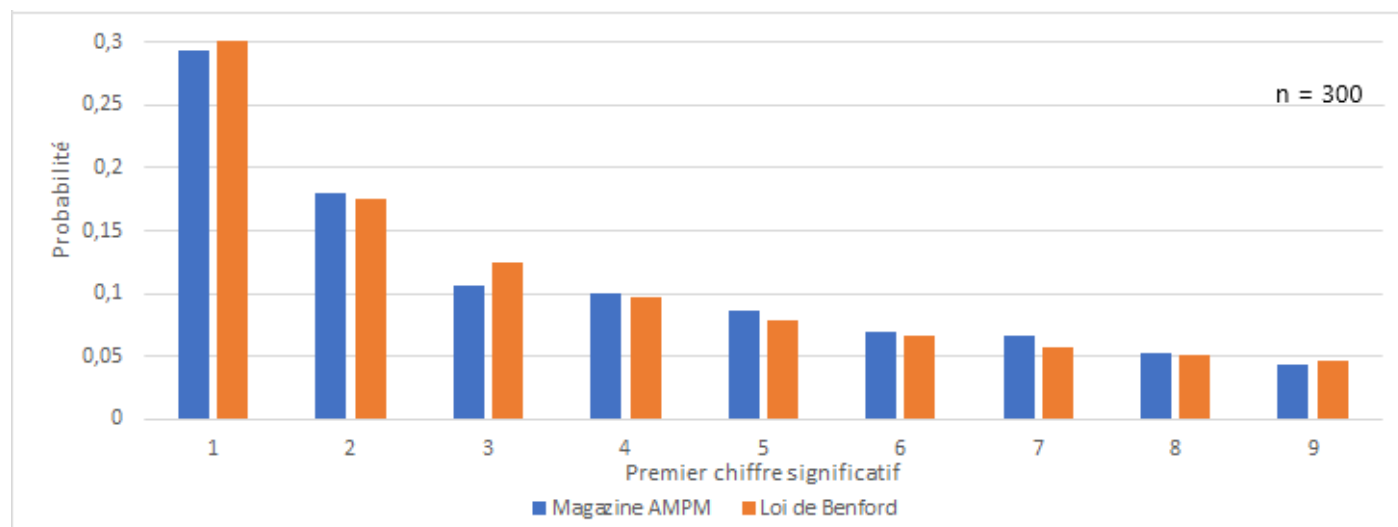


Figure 3 : Histogramme de la répartition du 1er chiffre significatif des prix du magazine AMPM en comparaison avec la loi de Benford

La répartition des chiffres significatifs des prix du magazine *AMPM* paraît fortement similaire à celle de la **loi de Newcomb-Benford**. Nous constatons tout de même une légère différence pour le chiffre 3.

Observons maintenant la répartition des données issues du journal *Les ECHOS*.

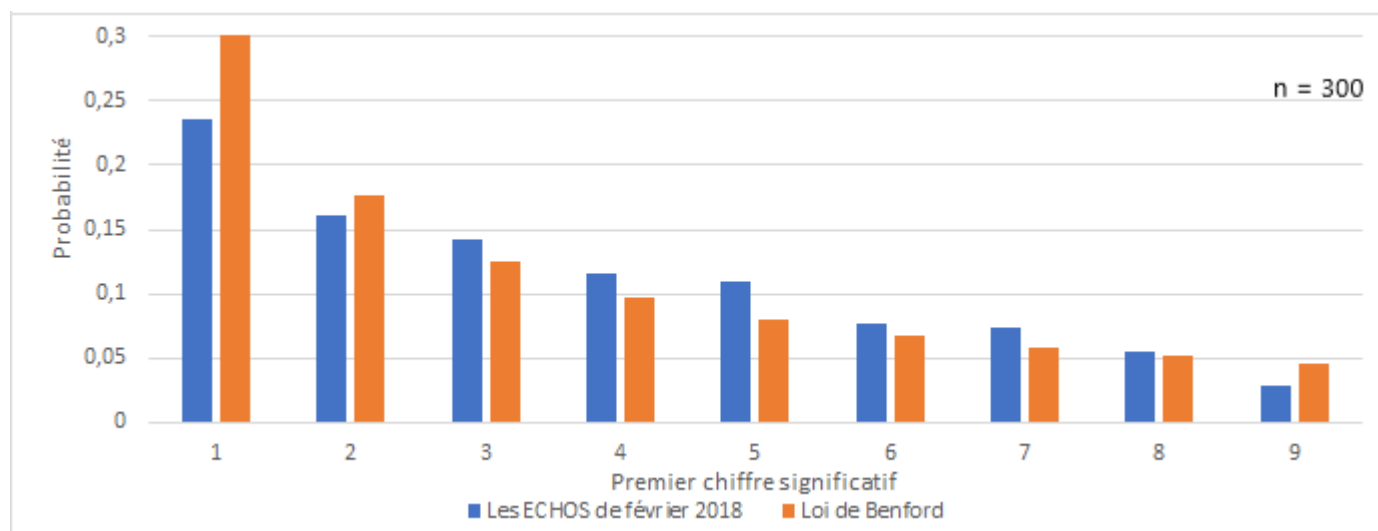


Figure 4 : Histogramme de la répartition du 1er chiffre significatif des nombres du journal LES ECHOS en comparaison avec la loi de Benford

Nous remarquons ici la même tendance décroissante. Cependant, les proportions des chiffres significatifs entre les données du journal et celles de la **loi de Newcomb-Benford** sont relativement différentes.

## Population des villes de France

Dans ce paragraphe, nous nous intéressons à la population des villes françaises. À l'aide des données de l'*INSEE*, nous répertorions environ 35000 premiers chiffres significatifs et regardons leur répartition.

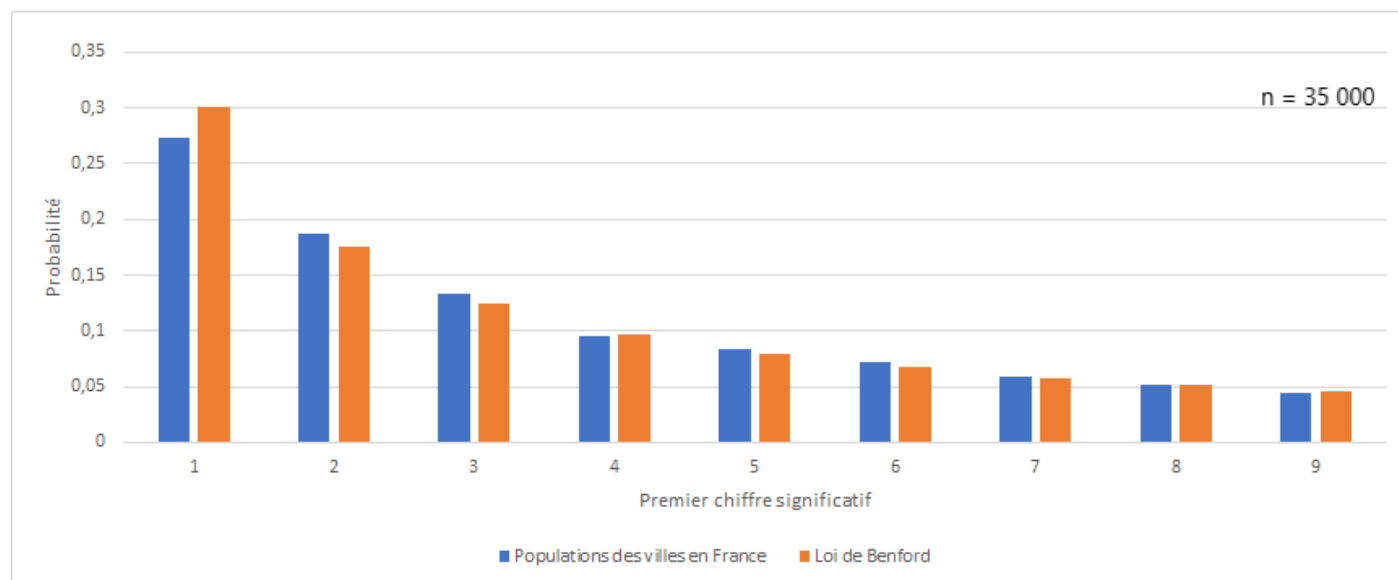


Figure 5 : Histogramme de la répartition du 1er chiffre significatif de la population française en comparaison avec la loi de Benford

Ici les répartitions sont fortement ressemblantes, c'est aussi le cas pour de nombreuses données démographiques naturelles. Nous aurions pu également analyser les codes postaux, la longueur des rivières ou encore la distance des villes de France à Paris.

## Passage journalier de vélos dans l'allée Beracasa à Montpellier

La ville de Montpellier étant en pleine transition écologique, elle ouvre de plus en plus l'accès aux vélos sur ses routes. Pour en mesurer l'impact, elle a mise en place des écompteurs dans plusieurs rues. Les données issues de ces compteurs sont en libre accès, nous nous sommes donc intéressés au nombre de passages journaliers de vélos dans l'allée Beracasa sur une année.

Nous obtenons la répartition suivante :

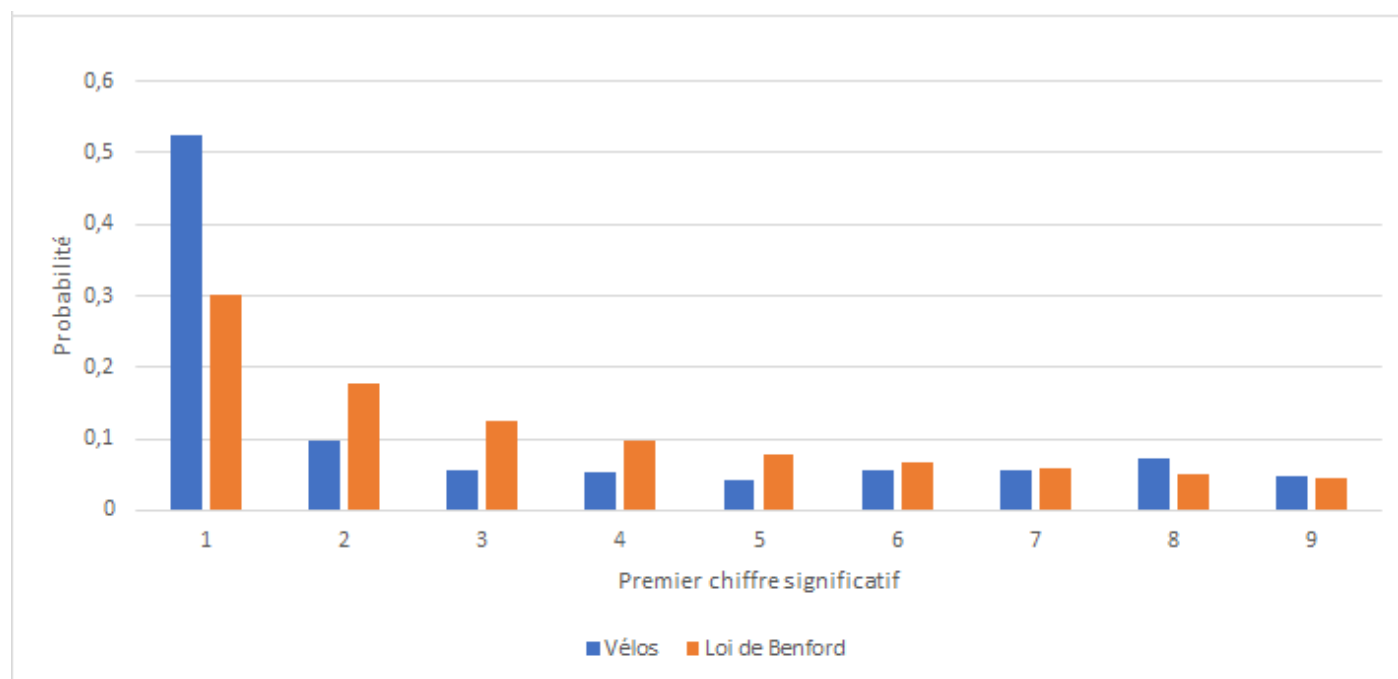


Figure 6 : Histogramme de la répartition du 1er chiffre significatif du passage journalier de vélos en comparaison avec la loi de Benford

Dans ce cas la proportion du chiffre 1 est de plus de 50% contre 30% pour la **loi de Newcomb-Benford**. La différence de répartition des chiffres 2, 3, 4, 5 est aussi notable, elle est même environ 2 fois moins élevée.

Visuellement, nous pourrions penser que la répartition de ces données ne suit pas la **loi de Newcomb-Benford**. Il est courant de ne pas retrouver la loi de Newcomb-Benford dans des données brutes comme celles-ci, on la retrouve empiriquement plus souvent dans des données dites de **deuxième génération** comme des sommes ou des produits. Ceci à été démontré par Jeff Boyle en 1994.

## Nombres générés par les humains

Dans ce paragraphe, nous tentons de reproduire à moindre échelle l'expérience de Theodore Preston Hill en 1988. Dans le cadre de son expérience le professeur Ted Hill demande à ses élèves (742) d'écrire un nombre de 6 chiffres au hasard sur un bout de papier, il recense ensuite le premier chiffre significatif de chacun de ces nombres dans le but de les comparer à la loi de Benford et à la répartition uniforme.

Notre expérience partageant le même objectif que celle de Hill, est basée sur un protocole légèrement différent. N'ayant pas une troupe d'élèves à disposition nous avons recueilli un total de 300 nombres.

Ces 300 nombres ont été obtenus de plusieurs manières, via des sondages sur internet ou sur les réseaux sociaux, en demandant directement à des personnes rencontrées au hasard, notre famille ou nos amis. Plus précisément, notre expérience a consisté à rechercher auprès de ces personnes un nombre à 2 chiffres, soit un nombre compris entre 10 et 99. Nous reviendrons plus loin sur l'importance que peut avoir ce détail.

D'après le biais d'équiprobabilité cité plus haut, si les nombres recensés pendant les expériences ont réellement été donnés de façon aléatoire la répartition du premier chiffre significatif devrait être comparable à une loi uniforme.

Comparons les répartitions obtenues durant les deux expériences :

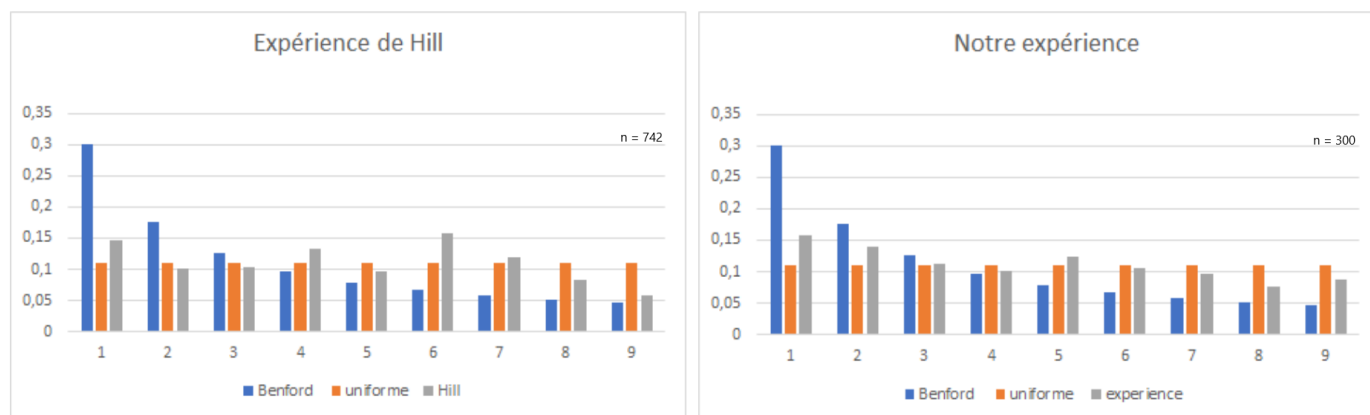


Figure 7 : Comparaison des histogrammes de la répartition du 1er chiffre significatif de l'expérience de Hill avec la loi de Benford et la loi uniforme puis de notre expérience avec ces deux mêmes lois

À première vue, dans les deux expériences la répartition du premier chiffre significatif ne semble pas suivre la loi de Newcomb-Benford (le chiffre 1 n'apparaît clairement pas aussi souvent par exemple), elle semble cependant plus proche de la loi uniforme sans tout de même y correspondre parfaitement.

Plusieurs facteurs pourraient expliquer les différences entre la distribution de la loi uniforme et la répartition du premier chiffre significatif de notre expérience, celui qui revient souvent est qu'un nombre donné au hasard par un humain est souvent influencé par son expérience, même inconsciemment. Par exemple sa date d'anniversaire, un évènement marquant ou son nombre préféré. Le fait d'avoir recueilli nos nombres par internet a aussi pu influencer le choix des personnes concernées. Un facteur psychologique est donc à prendre en considération pour approfondir la conclusion. Notons également que lors de notre expérience, la question posée stipulait de donner un nombre compris entre 10 et 99, soit un nombre à 2 chiffres. Ainsi, il est important de retenir que dès lorsqu'on demande un avec suffisamment de chiffres, plus le premier chiffre aura tendance à suivre la loi de Newcomb-Benford, mais la répartition des autres chiffres significatifs (deuxième, troisième, etc.) ne sera apparentée à aucune loi. De même si on demande un nombre avec peu de chiffres (choix d'un chiffre entre 1 et 9 par exemple) la répartition aura plutôt tendance à être uniforme. Ce phénomène a été exposé par A. Diekmann<sup>[n3]</sup> en 2007.

[n<sup>3</sup>] : Tiré de l'article *Not the First Digit! Using Benford's Law To Detect Fraudulent Scientific Data* écrit par A. Diekmann en 2007 (*Source*).

Après avoir observé ces quelques jeux de données, nous étions en mesure de dire si ces données semblaient ou non suivre la loi de Newcomb-Benford, le problème qui en découle est qu'une simple observation n'est pas très fiable, difficile de prendre une décision sur un constat visuel. En effet, se tromper dans

l'interprétation peut entraîner deux types d'erreur, la première étant de faussement déceler une fraude (ce que nous appellerons **le risque de première espèce**) et la deuxième de laisser passer une fraude. Ces erreurs ont un coût pour l'institut qui essaye de les réprimer, celui d'engager des démarches de détectations approfondies inutiles ou de ne pas percevoir les taxes dues dans le cas de la fraude fiscale par exemple.

Le but est donc de minimiser le coût que peuvent engendrer les erreurs susmentionnées, pour ce faire l'utilisation d'outils scientifiques est de rigueur. Les outils que nous aborderons dans la suite sont les tests d'adéquation, ces tests servent à vérifier si un ensemble de nombre suit ou non une loi de probabilité donnée (dans notre cas, c'est la loi de Newcomb-Benford).

## Tests

### Généralités sur les tests

#### METTRE SOURCE

Un test d'hypothèses (ou test statistique) est un procédé d'inférence statistique ayant pour but de fournir une règle de décision permettant ainsi, à partir de l'étude d'un ou plusieurs échantillons de données, d'indiquer si une hypothèse statistique concernant une population doit être acceptée ou rejetée.

Nous distinguons deux classes de tests :

- Les tests paramétriques sont l'étude de la moyenne, variance, ou de la fréquence des observations issues d'une distribution a priori paramétrée. Ils nécessitent un modèle à fortes contraintes (normalité des distributions ou approximation normale pour de grands échantillons). Ces hypothèses sont d'autant plus difficiles à vérifier que les effectifs étudiés sont plus réduits.
- Les tests non paramétriques sont l'étude des rangs des observations issues d'une distribution non paramétrée, mais quelconque. Ce sont des tests dont le modèle ne précise pas les conditions que doivent remplir les paramètres de la population dont a été extrait l'échantillon. Il n'y a donc pas d'hypothèse de normalité au préalable.

Lorsque les conditions nécessaires sont valides, les tests paramétriques sont plus puissants que les tests non paramétriques. Les tests non paramétriques s'utilisent dès lors que les conditions d'application des autres méthodes ne sont pas satisfaites, même après d'éventuelles transformations de variables et peuvent être employés pour des échantillons de taille très faible.

Comme nous l'avons précédemment énoncé, une inspection visuelle à elle seule ne permet pas d'affirmer ou infirmer si un jeu de données suit la loi de Newcomb-Benford. L'outil statistique permettant de le vérifier est le test d'adéquation à la loi de Newcomb-Benford.

Les tests d'adéquation servent à tester si un échantillon est distribué selon une loi de probabilité préalablement choisie. Ils permettent de décider, avec un certain seuil d'erreur, si les écarts présentés par l'échantillon par rapport aux valeurs théoriques sont dus au hasard, ou si au contraire ils sont significatifs.

## Hypothèse nulle et hypothèse alternative

Un test statistique étudie deux hypothèses opposées concernant une population : l'hypothèse nulle et l'hypothèse alternative.

L'hypothèse nulle, notée  $H_0$ , est l'hypothèse que l'on souhaite contrôler, elle repose sur le fait de dire qu'il n'existe pas de différence entre les paramètres comparés ou que la différence observée n'est pas significative et résulte des fluctuations d'échantillonnage.

À partir des échantillons de données, un test statistique permet de déterminer si on peut rejeter l'hypothèse nulle. La  $p$ -valeur sert de détermination. Si la  $p$ -valeur est inférieure au seuil de signification (appelé  $\alpha$ ), l'hypothèse nulle peut être rejetée.

L'hypothèse alternative, notée  $H_1$ , affirme qu'un paramètre de la population est plus petit, plus grand ou différent de la valeur hypothétique dans l'hypothèse nulle. Elle peut être vue comme la négation de  $H_0$  et est équivalente à dire que  $H_0$  est fausse. La décision de rejeter  $H_0$  signifie que  $H_1$  est réalisée ou que  $H_1$  est vraie.

On pense souvent à tort que les tests d'hypothèses statistiques visent à choisir l'hypothèse la plus probable parmi  $H_0$  et  $H_1$ . Néanmoins l'hypothèse nulle est formulée dans le but d'être rejetée. Le seuil de signification fixé est bas (généralement 0.05), et lorsque l'hypothèse nulle est rejetée cela prouve statistiquement que l'hypothèse alternative est vraie. En revanche, si l'hypothèse nulle ne peut être rejetée aucune preuve statistique ne montre que celle-ci est vraie. La raison est qu'il n'y a pas de valeur fixée assurant que la probabilité d'accepter à tort l'hypothèse nulle est petite.

Finalement, la décision d'accepter l'hypothèse nulle n'est pas équivalente à dire que  $H_0$  est vraie et que  $H_1$  est fausse, mais cela traduit uniquement l'idée selon laquelle il n'y a pas d'évidence nette pour que  $H_0$  soit fausse. Un test conclu donc à rejeter ou à ne pas rejeter l'hypothèse nulle mais jamais à l'accepter directement.

Dans la suite de notre projet, nous nous intéresserons aux données fiscales de 20 régions italiennes entre 2007 et 2011. Nous réaliserons donc pour chacune de ces régions et chacune des années le test d'hypothèses suivant :  $H_0$  : "La répartition du premier chiffre significatif suit la loi de Newcomb-Benford" contre  $H_1$  : "La répartition du premier chiffre significatif ne suit pas la loi de Newcomb-Benford".

## Les risques d'erreurs

Notons alors qu'aucun test d'hypothèses n'est fiable à 100%, un test étant basé sur des probabilités, il existe toujours un risque de tirer une mauvaise conclusion. Lorsqu'un test d'hypothèses est effectué, nous pouvons observer deux types d'erreurs, l'erreur de Type I dite erreur de première espèce et l'erreur de Type II dite erreur de seconde espèce. Les risques de ces deux erreurs sont inversement proportionnels et sont déterminés par le seuil de signification (ou région critique) et la puissance du test. Il est important de déterminer l'erreur qui présente les conséquences les plus graves dans notre cas avant de définir le risque que nous accepterons pour chaque erreur.

L'erreur de Type I consiste à rejeter l'hypothèse nulle alors que celle-ci est vraie. La probabilité de commettre une erreur de première espèce est représentée par  $\alpha$ , qui désigne le seuil de signification défini pour le test d'hypothèses. Ainsi, le seuil de signification du test s'énonce en probabilité :

$$\alpha = \mathbb{P}(\text{rejeter } H_0 | H_0 \text{ vraie}).$$

Décision d'après l'échantillon	Réalité sur la population	
	$H_0$ est vraie	$H_0$ est fausse
Ne pas rejeter $H_0$	Décision juste (probabilité = $1 - \alpha$ )	Erreur de seconde espèce : acceptation de $H_0$ alors que celle-ci est fausse (probabilité = $\beta$ )
Rejeter $H_0$	Erreur de première espèce : rejet de $H_0$ alors que celle-ci est vrai (probabilité = $\alpha$ )	Décision juste (probabilité = $1 - \beta$ )

Tableau 2 : Règle de décision et risques d'erreurs

Un niveau  $\alpha$  de 0.05 indique que nous sommes disposé à avoir 5% de chances de rejeter l'hypothèse nulle à tort. Pour réduire ce risque, il est possible d'utiliser une valeur  $\alpha$  plus faible. Cependant, cela implique d'être moins à même de détecter une vraie différence si celle-ci existe vraiment.

Dans notre contexte, l'erreur de Type I consiste à affirmer que les données ne suivent pas la loi de Newcomb-Benford alors que c'est le cas, soit faussement identifier une fraude.

L'erreur de Type II repose sur le fait de ne pas rejeter l'hypothèse nulle alors que celle-ci est fausse. La probabilité de commettre une erreur de seconde espèce est notée  $\beta$ , et dépend de la puissance du test. Il est possible de réduire le risque de deuxième espèce en faisant en sorte que le test soit suffisamment puissant. Pour cela, il est nécessaire que l'effectif d'échantillon soit suffisamment grand pour permettre la détection d'une différence réelle.

La probabilité de rejeter l'hypothèse nulle à tort vaut  $1 - \beta$ , il s'agit de la puissance du test.

Au sujet de la loi de Newcomb-Benford, notons que le test d'adéquation le plus populaire est le test du khi-deux de Pearson dont la puissance, associée au risque d'erreur de Type II est relativement faible. C'est pourquoi, d'autres tests ont été mis en place récemment. L'ensemble de ces tests seront développés par la suite.

G. Ducharme, S. Kaci et C. Vovor-Dassu<sup>3</sup> ont introduits de nouveaux tests d'adéquation pour cette loi, basés sur le principe des tests lisses. Ils ont également comparés ces tests aux meilleurs tests existants et ont montrés qu'ils seraient globalement plus performants. Notons aussi que la qualité d'un test dépend de sa puissance, plus forte elle est meilleure le test est. Cependant, différents tests peuvent conduire à des avantages différents. Par exemple, certains tests détecteront plus que d'autres la différence significative du premier chiffre significatif, quant à d'autres, ce sera un autre chiffre significatif.

3. Tests d'adéquations lisses pour la loi de Newcomb-Benford écrit en 2020 (*source*).



## Test du Khi-Deux

Les tests du  $\chi^2$  sont des tests d'hypothèses statistiques non paramétriques[<sup>n5</sup>]. Ceux-ci permettent de comparer la distribution observée dans un échantillon statistique avec une distribution théorique (*test d'ajustement*), à tester si deux caractères d'une population sont indépendants (*test d'indépendance*) et à tester si des échantillons sont issus d'une même population (*test d'homogénéité*). Le test lit l'écart critique dans la table de la loi du khi-deux.

Le déroulement du test se procède en 5 étapes :

- 1) On calcule les effectifs théoriques ( $n_{pj}$ ).
- 2) On calcule la valeur observée de la variable de test :

$$\chi^2 = \sum_{j=1}^k \frac{(n_j - n_{pj})^2}{n_{pj}}.$$

- 3) On cherche la valeur critique  $\chi_a^2$  dans la table de la loi du  $\chi^2$  à  $k - 1$  degrés de liberté.
- 4) Si  $\chi_a^2 < \chi^2$ , on ne rejette pas l'hypothèse  $H_0$  ("la distribution observée est conforme à la distribution théorique" avec un risque d'erreur  $\alpha$ ), sinon on la rejette.
- 5) Il faut vérifier que  $n_{pj} \geq 5$  pour tout  $j$ .

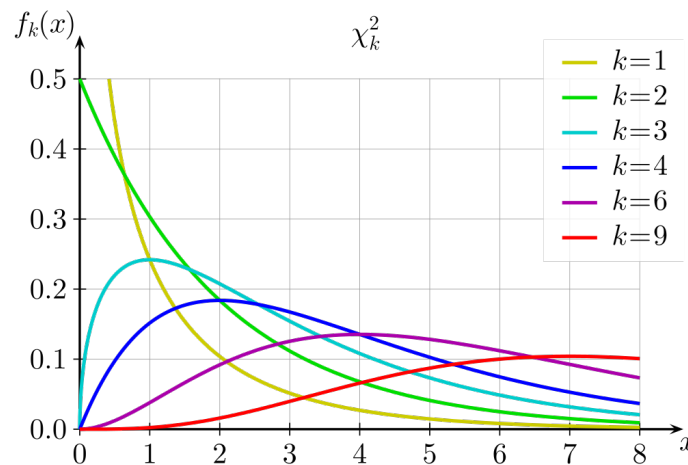


Figure 8 : Densité de la loi du Khi-Deux en fonction du nombre de degrés de liberté

[<sup>n5</sup>] : Cette partie a été fortement inspiré du site *Bibmath* (cf Bibliographie).

La loi du  $\chi^2$  à  $n$  degrés de liberté si elle est absolument continue, admet pour densité :

$$f(x) = \begin{cases} \frac{1}{2^{n/2}\Gamma(n/2)} e^{-x/2} x^{\frac{n}{2}-1} & \text{si } x > 0, \\ 0 & \text{sinon,} \end{cases}$$

où  $X$  admet alors pour espérance et variance  $E(X) = n$  et  $V(X) = 2n$ .

## Test de Freedman-Watson

Le test de Freedman peut être adaptée à des données discrètes, il permet de comparer la distribution du premier chiffre significatif d'un échantillon de données avec la distribution de la loi de Newcomb-Benford et ainsi affirmer si la répartition du premier chiffre significatif cet échantillon est bien conforme à la loi de Newcomb-Benford.

Spécifiquement, la statistique de test (dans le cas  $k = 1$ ) est donnée par :

$$U^2 = \frac{n}{9} \cdot \left[ \sum_{i=1}^8 \left( \sum_{j=1}^i (f_j^o - f_j^e) \right)^2 - \frac{1}{9} \cdot \left( \sum_{i=1}^8 \sum_{j=1}^i (f_i^o - f_i^e) \right)^2 \right],$$

Avec  $f_i^o$  la fréquence observée du chiffre  $i$  et  $f_i^e$  la fréquence attendue du chiffre  $i$ .

Notons que de plus grands écarts entre les fréquences conduisent à un plus grand  $U^2$ , ce qui rend le rejet plus probable. Ce test est reconnu comme plus performant que d'autres tests et a même été recommandé par la statisticienne M. Lesperance ainsi que ses collaborateurs (2016), puis également par Joenssen (2014).

## Tests lisses pour la Loi de Newcomb-Benford

La famille des tests lisses introduite par Neyman s'applique à des données autant discrètes que continues. Celle-ci est spécifique à la loi de probabilité sous  $H_0$ .

Il existe deux théorèmes essentiels tirés de l'article "*Smooth test of goodness-of-fit for directional and axial data*" écrit par BOULERICE B., DUCHARME G.R. en 1997 qui permettent de construire une famille de tests lisses pour l'hypothèse nulle  $H_0 : X \sim f(\cdot)$ . Ici  $f(\cdot)$  est la densité de la loi de Newcomb-Benford.

Le premier théorème nous dit :

**Théorème 1** : Soit  $X_1, \dots, X_n$  des copies indépendantes d'une variable aléatoire  $X$  de densité  $f(\cdot)$  par rapport à une mesure dominante  $\nu$ . Soit  $\{h_0(\cdot) \equiv 1, h_k(\cdot), k = 1, 2, \dots\}$  une suite de fonctions orthonormales par rapport à  $f(\cdot)$ ; plus précisément,  $\int h_k(x)h_{k'}(x)f(x)d\nu(x) = \delta_{kk'}$ , la fonction delta de Kronecker. Soit  $U_k = n^{-1/2} \sum_{i=1}^n h_k(X_i)$  et pour un entier  $K \geq 1$ , soit  $T_K = \sum_{k=1}^K U_k^2$ . Alors sous  $H_0$ ,  $T_K \xrightarrow{L} \chi_K^2$ , la loi khi-deux à  $K$  degrés de liberté, et un test de niveau asymptotique  $\alpha$  rejette  $H_0$  si la valeur observée de  $T_K$  dépasse  $x_{K,1-\alpha}^2$ , le quantile d'ordre  $1 - \alpha$  de cette loi  $\chi_K^2$ .

Nous avons donc nos statistiques de test  $T_K$ . Exprimons maintenant les  $h_k$ .

Dans la suite l'indice 0 dénote un opérateur probabiliste calculé sous  $H_0 : X$  suit  $f(\cdot)$ .

Nous avons également le théorème qui suit :

**Théorème 2** : Soit  $\mu_k = \mathbb{E}_0(X^k)$ ,  $k \geq 0$ . Soit aussi la matrice  $\mathbf{M}_k = [\mu_{i+i'}]_{i,i'=0,\dots,k-1}$ , le vecteur  $\boldsymbol{\mu}_k = (\mu_k, \mu_{k+1}, \dots, \mu_{2k-1})^T$  et la constante  $c_k = \mu_{2k} - \boldsymbol{\mu}_k^T \mathbf{M}_k^{-1} \boldsymbol{\mu}_k$ . Alors les polynômes

$$h_k(x) = c_k^{-1/2} \left( x^k - \left( 1, x, x^2, \dots, x^{k-1} \right) \mathbf{M}_k^{-1} \boldsymbol{\mu}_k \right)$$

satisfont la condition du Théorème 1.

D'après le Théorème 2, et Ducharme & Collab. nous avons pour  $0 < k \leq 5$  les  $h_k$  suivants :

$$\begin{aligned} h_1(x) &= -1.3979 + 0.4063x, \\ h_2(x) &= 2.2836 - 1.6128x + 0.18247x^2, \\ h_3(x) &= 4.0815 + 4.5719x - 1.2053x^2 + 0.0862x^3, \\ h_4(x) &= 8.0795 - 12.0946x + 5.1951x^2 - 0.8249x^3 + 0.0431x^4, \\ h_5(x) &= -18.1064 + 33.1385x - 19.7207x^2 + 5.0168x^3 - 0.5665x^4 + 0.0233x^5. \end{aligned}$$

Nous définissons :

$$\hat{K} = \arg \max_{1 \leq k \leq K_{\max}} \{T_k - k \log(n)\},$$

et la statistique de test  $T_{\hat{K}} \xrightarrow{L} \chi_1^2$  sous  $H_0$ .

Dans la suite, nous appliquerons les tests  $T_{\hat{K}}$  et  $T_2$  qui d'après Ducharme & Collab. sont les plus performants.

« « « < HEAD

# Application des tests

Nous traitons sur nos différents jeux de données et sur la fiscalité italienne de 20 régions, ces données nous serviront à appliquer différents tests (classique ou lisse). Chacun des tests teste l’hypothèse nulle “la répartition du premier chiffre significatif suit la loi de Newcomb-Benford” contre l’hypothèse alternative “la répartition du premier chiffre significatif ne suit pas la loi de Newcomb-Benford”. L’objectif ici est d’approuver nos résultats si les données étudiées suivent ou non la loi de Newcomb-Benford et de chercher à savoir si certaines régions d’Italie modifient ou non leurs chiffres, mais aussi de comparer les résultats de nos différents tests.

Nous retenons le test du khi-deux qui est le plus connu et utilisé, il semblerait cependant qu’il fasse partie des moins performants (Ducharme & Collab. 2020), le test de Freedman-Watson explicité plus haut. Pour ce qui est des tests lisses nous appliquerons les tests  $T_{\hat{K}}$  et  $T_2$  comme susmentionné.

## Application sur nos jeux de données

Dans cette partie nous cherchons à confirmer ou non nos hypothèses formulées à la vu des graphiques. Nous utilisons les test cités plus haut. Les données se présentent comme ci-après, chaque ligne correspond à une de nos expérience, les 9 premières colonnes (sans compter la première qui correspond à l’expérience) représente la répartition du premier chiffre significatif. Les 4 dernières colonnes quant à elles sont les p\_values des différents tests. Une case est rouge si l’hypothèse nulle est rejetée au seuil de 5% d’erreur.

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>Fibonacci</b>	0.301	0.177	0.125	0.095	0.080	0.067	0.056	0.053	0.045	<b>1.000</b>	<b>1.000</b>	<b>0.998</b>	<b>0.961</b>
<b>AMPM</b>	0.293	0.180	0.107	0.100	0.087	0.070	0.067	0.053	0.043	<b>0.990</b>	<b>0.939</b>	<b>0.872</b>	<b>0.623</b>
<b>Les echos</b>	0.235	0.161	0.142	0.116	0.110	0.077	0.074	0.055	0.029	<b>0.096</b>	<b>0.002</b>	<b>0.002</b>	<b>0.000</b>
<b>Population</b>	0.274	0.188	0.133	0.095	0.084	0.072	0.059	0.052	0.044	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>
<b>Velos</b>	0.525	0.098	0.056	0.053	0.042	0.056	0.056	0.071	0.048	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>
<b>Humains</b>	0.135	0.140	0.113	0.100	0.123	0.107	0.097	0.077	0.087	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>

Les tests sont unanimes, la suite de Fibonacci et les données relevées sur le magazine AMPM suivent la loi de Newcomb-Benford, du moins on ne rejette pas cette hypothèse, ce résultat est en accord avec nos suppositions.

Regardons maintenant les tests appliqués au jeu de données tirés du magazine *les échos*, nous avons remarqué une tendance décroissante similaire à celle de la répartition de la loi de Newcomb-Benford, avec cependant des valeurs relativement différentes. Les deux tests lisses ainsi que le test de Freedman donnent une p\_value proche de zéro qui nous mène à un fort rejet de l’hypothèse nulle (avec un risque d’erreur presque nul), contrairement au test du  $\chi^2$  (ici une p\_value d’environ 0,1 nous indique que si nous rejetons l’hypothèse nulle nous avons environ 10% de chance de nous tromper). On remarque ici que ce dernier est possiblement moins puissant que les autres.

Pour terminer tous les tests mènent à un rejet de l’hypothèse nulle “le jeu de donnée suit la loi de N-B” pour les jeux de données sur les population des villes de France, le passage journalier de vélos et les nombres générés par les humains. Pour les deux derniers pas de surprise, cependant le rejet de l’hypothèse nulle pour le jeu de données sur les population des villes de France nous paraît étonnant,

l'inspection visuelle nous donnait l'impression que la répartition du PCS collait à celle de la loi de Neewcomb-Benford, mais le nombre important de données à fortement amplifié les différences.

Les différents tests nous ont permis ici de déceler des différences avec la loi de N-B qui nous ont échappés à l'oeil nu, mais aussi de nous rendre compte de la différence d'efficacité entre les tests.

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.292	0.187	0.105	0.085	0.089	0.056	0.075	0.052	0.059	<b>0.716</b>	<b>0.135</b>	<b>0.396</b>	<b>0.332</b>
<b>2008</b>	0.295	0.170	0.111	0.089	0.098	0.059	0.062	0.056	0.059	<b>0.870</b>	<b>0.220</b>	<b>0.465</b>	<b>0.309</b>
<b>2009</b>	0.289	0.164	0.111	0.102	0.079	0.062	0.066	0.056	0.072	<b>0.637</b>	<b>0.146</b>	<b>0.103</b>	<b>0.091</b>
<b>2010</b>	0.298	0.154	0.121	0.092	0.085	0.069	0.066	0.052	0.062	<b>0.918</b>	<b>0.115</b>	<b>0.331</b>	<b>0.211</b>
<b>2011</b>	0.318	0.144	0.115	0.095	0.079	0.062	0.079	0.059	0.049	<b>0.746</b>	<b>0.248</b>	<b>0.443</b>	<b>0.380</b>

Tableau 3 : Tests sur la région d'Abruzzo

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.290	0.160	0.130	0.061	0.122	0.076	0.061	0.023	0.076	<b>0.267</b>	<b>0.289</b>	<b>0.708</b>	<b>0.428</b>
<b>2008</b>	0.282	0.168	0.107	0.107	0.084	0.092	0.031	0.061	0.069	<b>0.718</b>	<b>0.633</b>	<b>0.558</b>	<b>0.314</b>
<b>2009</b>	0.290	0.160	0.122	0.099	0.061	0.099	0.053	0.053	0.061	<b>0.894</b>	<b>0.225</b>	<b>0.647</b>	<b>0.388</b>
<b>2010</b>	0.290	0.168	0.115	0.092	0.076	0.122	0.023	0.053	0.061	<b>0.287</b>	<b>0.745</b>	<b>0.785</b>	<b>0.493</b>
<b>2011</b>	0.305	0.160	0.099	0.115	0.053	0.122	0.046	0.046	0.053	<b>0.348</b>	<b>0.505</b>	<b>0.874</b>	<b>0.611</b>

Tableau 4 : Tests sur la région de Basilicata

## Application a des données fiscales italiennes

Les données suivantes sont tirées Les données se présentent comme ci-après, chaque ligne correspond à une année, les 9 premières colonnes (sans compter la première qui correspond à l'année) représente la répartition du premier chiffre significatif. Les 4 dernières colonnes quant à elles sont les p\_values des différents tests que nous avons cité plus haut. Une case est rouge si l'hypothèse nulle est rejetée au seuil de 5% d'erreur.

Les 12 premiers tableaux ci-dessus ne montrent pas de suspicions de fraude dans leurs données avec les divers tests sélectionnés.

### Seulement un test

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.298	0.152	0.132	0.098	0.073	0.073	0.068	0.046	0.059	<b>0.816</b>	<b>0.328</b>	<b>0.481</b>	<b>0.270</b>
<b>2008</b>	0.301	0.164	0.127	0.093	0.073	0.068	0.073	0.042	0.059	<b>0.808</b>	<b>0.258</b>	<b>0.599</b>	<b>0.446</b>
<b>2009</b>	0.298	0.164	0.120	0.100	0.064	0.076	0.071	0.049	0.059	<b>0.764</b>	<b>0.481</b>	<b>0.378</b>	<b>0.270</b>
<b>2010</b>	0.308	0.166	0.115	0.103	0.061	0.081	0.073	0.051	0.042	<b>0.712</b>	<b>0.606</b>	<b>0.921</b>	<b>0.719</b>
<b>2011</b>	0.301	0.166	0.117	0.095	0.073	0.086	0.064	0.054	0.044	<b>0.933</b>	<b>0.379</b>	<b>0.814</b>	<b>0.521</b>

Tableau 5 : Tests sur la région de Calabria

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.304	0.193	0.119	0.095	0.061	0.058	0.053	0.056	0.061	<b>0.760</b>	<b>0.076</b>	<b>0.302</b>	<b>0.974</b>
<b>2008</b>	0.310	0.188	0.127	0.093	0.069	0.050	0.066	0.042	0.056	<b>0.823</b>	<b>0.177</b>	<b>0.566</b>	<b>0.654</b>
<b>2009</b>	0.315	0.193	0.116	0.087	0.082	0.048	0.061	0.040	0.058	<b>0.659</b>	<b>0.302</b>	<b>0.421</b>	<b>0.539</b>
<b>2010</b>	0.320	0.190	0.119	0.093	0.077	0.053	0.056	0.034	0.058	<b>0.690</b>	<b>0.338</b>	<b>0.370</b>	<b>0.332</b>
<b>2011</b>	0.312	0.198	0.127	0.079	0.077	0.061	0.045	0.048	0.053	<b>0.829</b>	<b>0.138</b>	<b>0.436</b>	<b>0.375</b>

Tableau 6 : Tests sur la région de la Lazio

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.279	0.210	0.142	0.078	0.055	0.059	0.068	0.064	0.046	<b>0.639</b>	<b>0.083</b>	<b>0.889</b>	<b>0.878</b>
<b>2008</b>	0.289	0.220	0.138	0.078	0.046	0.073	0.055	0.064	0.037	<b>0.439</b>	<b>0.347</b>	<b>0.840</b>	<b>0.602</b>
<b>2009</b>	0.280	0.220	0.142	0.078	0.041	0.069	0.064	0.064	0.041	<b>0.343</b>	<b>0.097</b>	<b>0.911</b>	<b>0.870</b>
<b>2010</b>	0.303	0.206	0.142	0.073	0.050	0.073	0.060	0.060	0.032	<b>0.590</b>	<b>0.091</b>	<b>0.731</b>	<b>0.441</b>
<b>2011</b>	0.284	0.211	0.151	0.069	0.055	0.073	0.041	0.073	0.041	<b>0.287</b>	<b>0.182</b>	<b>0.897</b>	<b>0.805</b>

Tableau 7 : Tests sur la région de Friuli-Venezia-Giulia

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.298	0.168	0.118	0.092	0.087	0.074	0.064	0.050	0.048	<b>0.725</b>	<b>0.449</b>	<b>0.479</b>	<b>0.225</b>
<b>2008</b>	0.297	0.175	0.115	0.089	0.085	0.073	0.070	0.045	0.051	<b>0.287</b>	<b>0.163</b>	<b>0.403</b>	<b>0.197</b>
<b>2009</b>	0.299	0.176	0.114	0.093	0.085	0.072	0.070	0.048	0.043	<b>0.500</b>	<b>0.395</b>	<b>0.759</b>	<b>0.526</b>
<b>2010</b>	0.296	0.172	0.119	0.096	0.082	0.071	0.069	0.054	0.043	<b>0.792</b>	<b>0.456</b>	<b>0.490</b>	<b>0.258</b>
<b>2011</b>	0.300	0.172	0.117	0.097	0.078	0.074	0.067	0.051	0.043	<b>0.775</b>	<b>0.258</b>	<b>0.783</b>	<b>0.507</b>

Tableau 8 : Tests sur la région de Lombardia

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.304	0.185	0.119	0.091	0.088	0.061	0.064	0.052	0.036	<b>0.616</b>	<b>0.430</b>	<b>0.626</b>	<b>0.459</b>
<b>2008</b>	0.300	0.186	0.112	0.096	0.085	0.065	0.065	0.044	0.047	<b>0.739</b>	<b>0.418</b>	<b>0.999</b>	<b>0.963</b>
<b>2009</b>	0.303	0.180	0.110	0.095	0.090	0.061	0.061	0.053	0.046	<b>0.766</b>	<b>0.373</b>	<b>0.925</b>	<b>0.787</b>
<b>2010</b>	0.304	0.185	0.111	0.094	0.086	0.055	0.055	0.054	0.041	<b>0.826</b>	<b>0.261</b>	<b>0.626</b>	<b>0.360</b>
<b>2011</b>	0.299	0.186	0.114	0.095	0.076	0.056	0.056	0.051	0.042	<b>0.680</b>	<b>0.311</b>	<b>0.621</b>	<b>0.384</b>

Tableau 9 : Tests sur la région de Piemonte



Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.306	0.194	0.116	0.089	0.089	0.043	0.054	0.047	0.062	<b>0.751</b>	<b>0.305</b>	<b>0.589</b>	<b>0.848</b>
<b>2008</b>	0.295	0.213	0.101	0.101	0.074	0.078	0.039	0.047	0.054	<b>0.648</b>	<b>0.907</b>	<b>0.903</b>	<b>0.750</b>
<b>2009</b>	0.291	0.217	0.101	0.097	0.074	0.078	0.035	0.062	0.047	<b>0.509</b>	<b>0.623</b>	<b>0.939</b>	<b>0.828</b>
<b>2010</b>	0.302	0.209	0.105	0.101	0.078	0.074	0.031	0.062	0.039	<b>0.558</b>	<b>0.565</b>	<b>0.770</b>	<b>0.470</b>
<b>2011</b>	0.310	0.202	0.112	0.101	0.070	0.074	0.035	0.054	0.043	<b>0.827</b>	<b>0.739</b>	<b>0.670</b>	<b>0.382</b>

Tableau 10 : Tests sur la région de Puglia

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.322	0.159	0.127	0.100	0.077	0.065	0.062	0.062	0.027	<b>0.788</b>	<b>0.329</b>	<b>0.801</b>	<b>0.562</b>
<b>2008</b>	0.300	0.168	0.126	0.099	0.081	0.048	0.075	0.060	0.042	<b>0.836</b>	<b>0.451</b>	<b>0.966</b>	<b>0.822</b>
<b>2009</b>	0.285	0.171	0.135	0.105	0.078	0.048	0.084	0.039	0.054	<b>0.436</b>	<b>0.991</b>	<b>0.889</b>	<b>0.655</b>
<b>2010</b>	0.300	0.168	0.120	0.114	0.078	0.057	0.072	0.048	0.042	<b>0.935</b>	<b>0.534</b>	<b>0.922</b>	<b>0.968</b>
<b>2011</b>	0.294	0.156	0.117	0.120	0.087	0.048	0.072	0.048	0.057	<b>0.538</b>	<b>0.861</b>	<b>0.736</b>	<b>0.437</b>

Tableau 11 : Tests sur la région de Trentino-alto

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.270	0.230	0.095	0.081	0.122	0.054	0.081	0.041	0.027	<b>0.708</b>	<b>0.633</b>	<b>0.728</b>	<b>0.902</b>
<b>2008</b>	0.270	0.243	0.108	0.068	0.108	0.054	0.081	0.014	0.054	<b>0.563</b>	<b>0.806</b>	<b>0.939</b>	<b>0.825</b>
<b>2009</b>	0.284	0.203	0.122	0.054	0.149	0.054	0.068	0.027	0.041	<b>0.486</b>	<b>0.606</b>	<b>0.799</b>	<b>0.870</b>
<b>2010</b>	0.284	0.216	0.108	0.027	0.162	0.068	0.054	0.041	0.041	<b>0.185</b>	<b>0.645</b>	<b>0.923</b>	<b>0.983</b>
<b>2011</b>	0.284	0.176	0.135	0.041	0.135	0.081	0.054	0.027	0.068	<b>0.501</b>	<b>0.984</b>	<b>0.908</b>	<b>0.654</b>

Tableau 12 : Tests sur la région de Valle D'aosta

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.287	0.195	0.123	0.077	0.085	0.079	0.051	0.049	0.054	<b>0.798</b>	<b>0.235</b>	<b>0.904</b>	<b>0.706</b>
<b>2008</b>	0.272	0.205	0.123	0.077	0.082	0.069	0.067	0.044	0.062	<b>0.460</b>	<b>0.088</b>	<b>0.602</b>	<b>0.362</b>
<b>2009</b>	0.279	0.203	0.123	0.077	0.079	0.072	0.064	0.051	0.051	<b>0.817</b>	<b>0.035</b>	<b>0.830</b>	<b>0.574</b>
<b>2010</b>	0.282	0.197	0.131	0.074	0.085	0.056	0.079	0.038	0.056	<b>0.285</b>	<b>0.389</b>	<b>0.845</b>	<b>0.631</b>
<b>2011</b>	0.290	0.195	0.126	0.079	0.077	0.067	0.069	0.038	0.059	<b>0.674</b>	<b>0.172</b>	<b>0.793</b>	<b>0.737</b>

Tableau 13 : Tests sur la région de la Sicilia

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.310	0.165	0.136	0.117	0.057	0.057	0.062	0.040	0.057	<b>0.184</b>	<b>0.172</b>	<b>0.770</b>	<b>0.714</b>
<b>2008</b>	0.313	0.167	0.139	0.107	0.062	0.062	0.065	0.041	0.043	<b>0.619</b>	<b>0.268</b>	<b>0.666</b>	<b>0.374</b>
<b>2009</b>	0.320	0.170	0.138	0.105	0.062	0.064	0.062	0.043	0.036	<b>0.613</b>	<b>0.408</b>	<b>0.316</b>	<b>0.139</b>
<b>2010</b>	0.325	0.172	0.138	0.103	0.067	0.062	0.055	0.053	0.024	<b>0.297</b>	<b>0.527</b>	<b>0.094</b>	<b>0.042</b>
<b>2011</b>	0.322	0.170	0.134	0.105	0.071	0.057	0.060	0.048	0.033	<b>0.704</b>	<b>0.763</b>	<b>0.270</b>	<b>0.118</b>

Tableau 15 : Tests sur la région de Veneto

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.309	0.160	0.096	0.094	0.067	0.102	0.051	0.054	0.067	<b>0.006</b>	<b>0.201</b>	<b>0.026</b>	<b>0.063</b>
<b>2008</b>	0.314	0.160	0.093	0.100	0.065	0.100	0.044	0.062	0.064	<b>0.003</b>	<b>0.404</b>	<b>0.034</b>	<b>0.114</b>
<b>2009</b>	0.310	0.176	0.091	0.091	0.071	0.093	0.051	0.060	0.058	<b>0.068</b>	<b>0.090</b>	<b>0.096</b>	<b>0.236</b>
<b>2010</b>	0.310	0.172	0.091	0.087	0.082	0.087	0.049	0.065	0.056	<b>0.094</b>	<b>0.069</b>	<b>0.089</b>	<b>0.186</b>
<b>2011</b>	0.310	0.172	0.087	0.083	0.087	0.078	0.060	0.058	0.064	<b>0.101</b>	<b>0.007</b>	<b>0.024</b>	<b>0.106</b>

Tableau 15 : Tests sur la région de Campania

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.296	0.201	0.103	0.078	0.069	0.069	0.057	0.066	0.060	<b>0.482</b>	<b>0.010</b>	<b>0.132</b>	<b>0.398</b>
<b>2008</b>	0.310	0.204	0.101	0.080	0.066	0.066	0.060	0.063	0.049	<b>0.634</b>	<b>0.035</b>	<b>0.311</b>	<b>0.979</b>
<b>2009</b>	0.305	0.201	0.103	0.072	0.075	0.060	0.060	0.060	0.063	<b>0.430</b>	<b>0.005</b>	<b>0.090</b>	<b>0.574</b>
<b>2010</b>	0.290	0.201	0.112	0.072	0.072	0.055	0.066	0.063	0.069	<b>0.225</b>	<b>0.019</b>	<b>0.044</b>	<b>0.250</b>
<b>2011</b>	0.299	0.195	0.115	0.072	0.060	0.069	0.063	0.072	0.055	<b>0.384</b>	<b>0.019</b>	<b>0.128</b>	<b>0.423</b>

Tableau 16 : Tests sur la région d'Emilia-romagna

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.272	0.209	0.085	0.089	0.077	0.098	0.043	0.047	0.081	<b>0.044</b>	<b>0.078</b>	<b>0.172</b>	<b>0.148</b>
<b>2008</b>	0.268	0.213	0.089	0.081	0.055	0.111	0.055	0.034	0.094	<b>0.001</b>	<b>0.597</b>	<b>0.044</b>	<b>0.073</b>
<b>2009</b>	0.294	0.209	0.077	0.098	0.060	0.094	0.060	0.034	0.077	<b>0.047</b>	<b>0.106</b>	<b>0.267</b>	<b>0.418</b>
<b>2010</b>	0.289	0.196	0.094	0.102	0.047	0.094	0.072	0.030	0.077	<b>0.043</b>	<b>0.570</b>	<b>0.285</b>	<b>0.320</b>
<b>2011</b>	0.306	0.191	0.089	0.115	0.047	0.089	0.060	0.043	0.060	<b>0.286</b>	<b>0.246</b>	<b>0.650</b>	<b>0.800</b>

Tableau 17 : Tests sur la région de la Liguria

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.268	0.172	0.092	0.121	0.084	0.063	0.054	0.071	0.075	<b>0.199</b>	<b>0.108</b>	<b>0.049</b>	<b>0.026</b>
<b>2008</b>	0.285	0.172	0.092	0.113	0.084	0.071	0.059	0.054	0.071	<b>0.593</b>	<b>0.163</b>	<b>0.208</b>	<b>0.122</b>
<b>2009</b>	0.280	0.180	0.088	0.117	0.088	0.071	0.067	0.042	0.067	<b>0.497</b>	<b>0.139</b>	<b>0.420</b>	<b>0.200</b>
<b>2010</b>	0.285	0.180	0.096	0.096	0.109	0.054	0.071	0.038	0.071	<b>0.269</b>	<b>0.473</b>	<b>0.429</b>	<b>0.239</b>
<b>2011</b>	0.293	0.197	0.088	0.096	0.105	0.042	0.071	0.054	0.054	<b>0.376</b>	<b>0.501</b>	<b>0.716</b>	<b>0.567</b>

Tableau 18 : Tests sur la région de Marche

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.243	0.125	0.191	0.118	0.074	0.074	0.074	0.051	0.051	<b>0.284</b>	<b>0.175</b>	<b>0.165</b>	<b>0.162</b>
<b>2008</b>	0.221	0.110	0.213	0.118	0.059	0.103	0.044	0.074	0.059	<b>0.007</b>	<b>0.923</b>	<b>0.047</b>	<b>0.033</b>
<b>2009</b>	0.250	0.110	0.184	0.110	0.081	0.103	0.051	0.066	0.044	<b>0.156</b>	<b>0.860</b>	<b>0.110</b>	<b>0.116</b>
<b>2010</b>	0.228	0.132	0.169	0.118	0.096	0.081	0.066	0.037	0.074	<b>0.233</b>	<b>0.543</b>	<b>0.079</b>	<b>0.051</b>
<b>2011</b>	0.235	0.140	0.169	0.096	0.103	0.081	0.081	0.029	0.066	<b>0.252</b>	<b>0.219</b>	<b>0.117</b>	<b>0.087</b>

Tableau 19 : Tests sur la région de Molise

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.310	0.167	0.093	0.053	0.050	0.119	0.064	0.109	0.034	<b>0.000</b>	<b>0.761</b>	<b>0.010</b>	<b>0.000</b>
<b>2008</b>	0.310	0.172	0.095	0.106	0.069	0.093	0.077	0.024	0.053	<b>0.048</b>	<b>0.659</b>	<b>0.956</b>	<b>0.772</b>
<b>2009</b>	0.316	0.175	0.101	0.095	0.085	0.069	0.093	0.032	0.034	<b>0.084</b>	<b>0.476</b>	<b>0.825</b>	<b>0.741</b>
<b>2010</b>	0.318	0.175	0.095	0.090	0.093	0.061	0.095	0.034	0.037	<b>0.042</b>	<b>0.129</b>	<b>0.962</b>	<b>0.816</b>
<b>2011</b>	0.314	0.174	0.103	0.087	0.092	0.066	0.103	0.029	0.032	<b>0.006</b>	<b>0.832</b>	<b>0.860</b>	<b>0.880</b>

Tableau 20 : Tests sur la région de la Sardegna

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.331	0.188	0.125	0.059	0.066	0.070	0.035	0.059	0.066	<b>0.171</b>	<b>0.112</b>	<b>0.045</b>	<b>0.013</b>
<b>2008</b>	0.328	0.195	0.105	0.080	0.070	0.059	0.042	0.059	0.063	<b>0.526</b>	<b>0.169</b>	<b>0.086</b>	<b>0.713</b>
<b>2009</b>	0.321	0.206	0.094	0.084	0.073	0.056	0.038	0.066	0.063	<b>0.255</b>	<b>0.087</b>	<b>0.080</b>	<b>0.841</b>
<b>2010</b>	0.341	0.195	0.101	0.084	0.066	0.059	0.042	0.045	0.066	<b>0.346</b>	<b>0.184</b>	<b>0.043</b>	<b>0.358</b>
<b>2011</b>	0.369	0.185	0.101	0.091	0.073	0.042	0.056	0.038	0.045	<b>0.280</b>	<b>0.563</b>	<b>0.023</b>	<b>0.030</b>

Tableau 21 : Tests sur la région de Toscana

Année	Chiffre significatif									p-value			
	1	2	3	4	5	6	7	8	9	$\chi^2$	<i>Freedman</i>	$T_2$	$T_{\hat{K}}$
<b>2007</b>	0.38	0.207	0.130	0.065	0.109	0.054	0.011	0.022	0.022	<b>0.210</b>	<b>0.786</b>	<b>0.030</b>	<b>0.021</b>
<b>2008</b>	0.37	0.207	0.130	0.065	0.120	0.043	0.022	0.022	0.022	<b>0.241</b>	<b>0.426</b>	<b>0.052</b>	<b>0.025</b>
<b>2009</b>	0.37	0.207	0.130	0.076	0.120	0.043	0.022	0.022	0.011	<b>0.196</b>	<b>0.701</b>	<b>0.022</b>	<b>0.017</b>
<b>2010</b>	0.37	0.217	0.120	0.076	0.109	0.065	0.011	0.000	0.033	<b>0.136</b>	<b>0.871</b>	<b>0.035</b>	<b>0.019</b>
<b>2011</b>	0.38	0.217	0.109	0.087	0.109	0.054	0.011	0.011	0.022	<b>0.154</b>	<b>0.683</b>	<b>0.018</b>	<b>0.013</b>

Tableau 22 : Tests sur la région de Umbria

## CAS COVID

Un article écrit par Junyi Zhang qui teste le nombre de cas du Covid-19 en Chine avec la loi de Newcomb-Benford. Il trouve une p-value de 92.8 % en faveur que les numéros respectent la loi de Newcomb-Benford. Pour confirmer le résultat, Ducharme et ses coloborateurs ont utilisée les mêmes données pour retrouver la même chose que le chercheur chinois malheureusement cela ne collait pas avec les résultats attendus. Nous-mêmes nous sommes allées chercher les données de cette article, mais le lien des données n'existe plus sur Wikipédia. On peut donc penser que les données ne sont pas cohérentes ou qu'il y a eu une erreur dans son étude ou frauder sur les chiffres.

## Conclusion

(Rappel de la problématique) Ce projet nous a permis d'approfondir nos connaissances, de pouvoir faire des expériences nous-mêmes et de nous initier au vaste problème pour détecter une fraude. Et plus particulièrement à la loi de Benford-Newcomb en exposant sa genèse mais aussi de l'étudier avec divers tests. (Reponse) Dans le déroulement de notre projet, nous avons appris que la loi de Benford-Newcomb n'est pas applicable à toutes les données comme l'expérience de Hill citée plus haut qui sont influencées par la pensée de l'Homme. De plus, lors des applications des tests, nous avons remarqué que certains tests avaient une p-value inférieure à 5%. Les tests paraissent fiables mais pour " détecter une fraude plus y a de test mieux c'est " Car les tests ne vont pas dans les mêmes directions car chaque tests mettent en valeur leurs "propriétés" c'est- à-dire que un permet de ... et un autre de .. (Ouverture) Pour avoir un autre point de vue, il serait bien d'appliquer la loi de Newcomb-Benford avec le deuxième chiffre significatif et lui appliquer les différents tests pour permettre à une autre conclusion ?

# Bibliographie

## Génèse de la loi :

V. GENEST, C. GENEST *La loi de Newcomb-Benford ou la loi du premier chiffre significatif* (2011) (*Source*)

WIKIPEDIA L'ENCYCLOPEDIE LIBRE *Loi de Benford* (dernière mise à jour 2021) (*Source*)

T. P. HILL *Random-number guessing and the first digit phenomenon* (1988) (*Source*)

## Test d'hypothèses :

G. R. DUCHARME, S. KACI, C. VOVOR-DASSU *Tests d'adéquations lisses pour la loi de Newcomb-Benford* (2020) (*Source*)

LENOIR *Les tests d'hypothèses* (*Source*)

MINITAB *Tests d'hypothèses* (*Source*)

WIKIPEDIA L'ENCYCLOPEDIE LIBRE *Test statistique* (dernière mise à jour 2021) (*Source*)

J. J. RUCH *Statistique : Tests d'hypothèses* (2012/2013) (*Source*)

## Khi-deux :

BIBMATH *Loi du Khi-deux* (*Source*)

WIKIPEDIA L'ENCYCLOPEDIE LIBRE *Test du Khi-deux de Pearson* (dernière mise à jour 2021) (*Source*)

BIBMATH *Tests du Khi-deux* (*Source*)

WIKIMEDIA COMMONS *Khi-deux* (*Source*)

## Freedman :

D. JOENSSEN, T. MUELLERLEILE *Package R "BenfordTest"* (2015) (*Source*)

D. JOENSSEN *Testing for Benford's Law : A Monte Carlo Comparison of Methods* (2014) (*Source*)

LESPERANCE M., REED W.J., STEPHENS M.A., TSAO C., WILTON B., *Assessing conformance with Benford's law : Goodness-of-fit tests and simultaneous confidence interval* (2016) (*Source*)

M. AUSLOOS, R. CERQUETI, T. A. MIR *Data science for assessing possible tax income manipulation : the case of Italy* (2017) (*Source*)

## Test lisse :

B. BOULERICE, G.R. DUCHARME *Smooth test of goodness-of-fit for directional and axial data*, Journal Multivariate Analysis 60 : 154-175 (1997) (*Source*)

G.R. DUCHARME, S. KACI, C. VOVOR-DASSU *article Tests d'adéquations lisse pour la loi de Newcomb-Benford* (2020) (*Source*)

G.R. DUCHARME, S. KACI, C. VOVOR-DASSU *BenfordSmoothTest* (2020) (*Source*)

NEYMAN J. *Smooth tests for goodness-of-fit*, Skand. Aktuarietidskr 20 : 150 – 199 (1937)

**Cas fiscalité italienne :**

*Data science for assessing possible tax incomemanipulation : the case of Italy* (2017)

**Cas Covid-19 :**

JUNYI ZHAN (2020) *Testing Case Number of Coronavirus Disease 2019 in China with Newcomb-Benford Law* (2020) (*Source*)